1  **NGSpop: A desktop software that supports population studies by**

2  **identifying sequence variations from next-generation sequencing data**

3  Short title: NGSpop, a desktop software for identifying sequence variations from next-

4  generation sequencing data

5  Dong-Jun Lee[1,*], Taesoo Kwon[2], Hye-Jin Lee[1], Yun-Ho Oh[1], Jin-Hyun Kim[1], Tae-Ho Lee[1]

6

7  [1]Genomics Division, National Institute of Agricultural Science, Jeonju, Republic of Korea

8  [2]Corporate R&D Center, Cloud9, Cheongju-si, Republic of Korea

9  * Corresponding author

10  E-mail: leemoses1004@gmail.com (D-JL)

11

12

# Abstract

14  Next-generation sequencing (NGS) is widely used in all areas of genetic research, such

15  as for genetic disease diagnosis and breeding, and it can produce massive amounts of

16  data. The identification of sequence variants is an important step when processing large

17  NGS datasets; however, currently, the process is complicated, repetitive, and requires

18  concentration, which can be taxing on the researcher. Therefore, to support researchers

19  who are not familiar with bioinformatics in identifying sequence variations regularly from

20  large datasets, we have developed a fully automated desktop software, NGSpop. NGSpop

21  includes functionalities for all the variant calling and visualization procedures used when

22  processing NGS data, such as quality control, mapping, filtering details, and variant

23  calling. In the variant calling step, the user can select the GATK or DeepVariant algorithm

24    for variant calling. These algorithms can be executed using pre-set pipelines and options

25    or customized with the user-specified options. NGSpop is implemented using JavaFX

26    (version 1.8) and can thus be run on Unix like operating systems such as Ubuntu Linux

27    (version 16.04, 18.0.4). Although there are several pipelines and visualization tools

28    available for NGS data analysis, most integrated environments do not support batch

29    processes; thus, variant detection cannot be automated for population-level studies. The

30    NGSpop software, developed in this study, has an easy-to-use interface and helps in rapid

31    analysis of multiple NGS data from population studies.

32

33

34

## Introduction

36    Next-generation sequencing (NGS) is widely used in all areas of genetic research, such

37    disease diagnosis and breeding, this is in part because it is a useful tool for the detection

38    of sequence variations [1-3]. NGS technology was originally used to study individuals

39    and small samples, but more recently, it has been used to study cohort-level populations.

40    In a medical study, such as that by the Undiagnosed Diseases Network (UDN) showed

41    that a genetic diagnosis with NGS is valid, even if the disease is undiagnosed [4].

42    According to NGS, 21% were changed in therapy, 37% in diagnostic testing, and 36% in

43    variant-specific genetic counseling. NGS has also been used to construct an ultra-high-

44    density genetic map for the identification of molecular markers for agricultural research

45    [5,6]. The research showed that a genetic breeding with NGS is a valid and reliable tool

46    to develop useful characters. NGS produces a large amount of data, especially for studies

47    involving genetic diseases and breeding at the population level. The identification of

48    sequence variants in these large datasets is one of the most important processing steps;

49    however, currently, sequence variation detection is both complicated and repetitive.

50    Genomics consortia, such as the 1000 genome project [7], provide shell scripts that

51    implement a standard operation procedure (SOP) for variant detection, which helps to

52    standardize the process (https://github.com/ekg/1000G-integration). However, most of

53    the SOP shell scripts in use are difficult to understand and automate. There are several

54    workflows and tools available that include quality control (QC), mapping and the calling,

55    annotation, and visualization of variations. Some tools have too many functions, and

56    consequently, they can be difficult to learn and often require official training. Furthermore,

57    for some tools, the lack of tool integration, and the many options included in their

58    functionality, can confuse the user and considering the available options can be time

59    consuming. Many pipelines and workflows have been developed by commercial and

60    open-source communities to support NGS data analysis. Pipelines such as the

61    ngs_backbone [8] and GATK [9] provide simple commands to perform a complete NGS

62    data analysis. Most pipelines offer only a command-line interface, and thus the user needs

63    to be trained in Unix/Linux commands, shell scripts, or Python. It is difficult to automate

64    variant detection in population-level studies. Galaxy [10] and the CLC genomics

65    workbench [11] provide users with easy-to-use graphical user interfaces (GUIs).

66    Although there are many pipelines and integrated environments for NGS data analysis,

67    each has its own strengths and limitations (Table 1).

68

3

69 **Table 1. Comparison of the user-friendly graphic interfaces and functions of the**

70 **SNP analysis pipelines**.

| Analysis \ Name | Annovar | Ngs_backbone | inGAP | Galaxy | CLC genomics workbench | NGSpop |
|---|---|---|---|---|---|---|
| Quality Control (QC) | | O | O | O | O | O |
| Read Mapping | | O | O | O | O | O |
| Variant calling (GATK) | | O | O | O | O | O |
| Variant calling (DeepVariant) | | | | | | O |
| Variant annotation | O | | | O | O | O |
| Visualization | | | | | | O |
| Manual mode | O | O | O | | O | O |
| Batch mode | | | | O | O | O |

71

72

73 To support sequence variation detection in population-level genomics studies, we have

74 developed a desktop software, NGSpop. The software accepts multiple NGS datasets and

75 allows the user to select between the GATK or DeepVariant [12] calling algorithms. The

76 functionalities for variant detection include QC, mapping, filtering, variant calling, and

77 visualization. Moreover, NGSpop has two modes of action: a one-step mode that supports

78 batch identification of variants and a step-by-step mode in which the user can verify the

79 result of each step. When the user selects the one-step mode, NGSpop can be executed

80 using pre-set options to exclude the time-consuming steps. NGSpop can only be used

81 with Linux operating systems.

82

83 # Implementation

84    NGSpop was implemented using JavaFX (version 1.8), and the tools employed within it

85    were compiled on Ubuntu Linux (version 18.0.4). The GNU compiler collection version

86    7.2.0, for Ubuntu Linux, was used as a C-language compiler.

87

## 88    Tools used in the pipeline

89    The tools included in NGSpop were carefully chosen according to the pipeline of the

90    National Agricultural Biotechnology Information Center (NABIC, Republic of Korea;

91    Fig. 1). NGS data need to be evaluated for QC, and for this purpose, NGSpop includes

92    FastQC (version 0.11.5). Filtering and trimming of the NGS data is mandatory, depending

93    on the sequence quality, and for this step, NGSpop employs TrimmOmatic (version 0.36)

94    [13]. After the QC step, sequence reads can be mapped in NGSpop against a reference

95    genome using an alignment tool, such as BWA (version 0.7.16a) [14], and SAMtools [15]

96    is used for file format conversion and indexing. Mate-pair information cannot be

97    concordant with the sample library information and should be fixed. If sequence reads

98    can be mapped to more than two loci, then the duplicate reads should be removed, and

99    Picard (version 2.9.4) is used for this in NGSpop. For SNP/INDEL identification, the user

100   can select SNP/INDEL identification algorithms from the Genome Analysis Toolkit

101   (version 3.7.0) or DeepVariant (version 0.5.1). Currently, DeepVariant is only supported

102   by the Linux operating system, and consequently, this system is required to run NGSpop.

103   To annotate the identified SNP/INDELs, SnpEff is used (version 4.3q) [16]. The

104   identified and annotated variants are visualized using JBrowser software (version 1.12.3)

105   [17]. All the tools integrated into NGSpop are summarized in Table 2.

106     **Fig 1. NGS data analysis pipeline used in the NGSpop software.** The variant analysis

107     protocol and tools are chosen according to the pipeline of the National Agricultural

108     Biotechnology Information Center (NABIC, Republic of Korea).

109

110     **Table 2. Tools included in the NGSpop software**

| Step | Tool | Version | Reference |
|------|------|---------|-----------|
| QC | FastQC | 0.11.5 | (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/) |
| | TrimmOmatic | 0.36 | [13] |
| Alignment | BWA | 0.7.16a | [14] |
| Post-processing | Samtools | 0.1.18 | [15] |
| | Picard | 2.9.4 | (http://broadinstitute.github.io/picard/) |
| | BamTools | 2.4.2 | [21] |
| | GATK (IndelRealigner) | 3.7.0 | [9] |
| Variant calling | GATK (HaplotypeCaller) | 3.7.0 | [9] |
| | GATK (UnifiedGenotyper) | 3.7.0 | [9] |
| | DeepVariant | 0.5.1 | [12] |
| Variant annotation | SnpEff | 4.3q | [16] |
| Visualization | Jbrowser | 1.12.3 | [17] |

111     The tools are listed in the order of their use in the pipeline.

112

# Project creation and importing input files

114     The user must create a project and specify the data files, including fastq files of

115     sequencing reads and a reference file in the FASTA format (Suppl. 2a). Fastq files of the

116     sequencing reads can be multiple pairs of forward and reverse reads for population studies.

117     Only fastq files produced by the Illumina platform can be processed using NGSpop. For

118     the convenience of users, NGSpop can download a reference file from a genomic database

119     such as the NCBI, Ensemble, and the NABIC server through the application program

6

120    interface. When the index file of the reference sequence does not exist, NGSpop performs

121    indexing of the reference file.

122

## Step-by-step mode

124    NGSpop provides the user with a step-by-step mode, in which they can investigate each

125    step of the analysis. The user can change or execute each option during each step, and

126    changes will become the default options for the same step in each subsequent run (Suppl.

127    2b). To monitor the progress of each step, NGSpop provides the user with a log window.

128

## One-step mode

130    To automate NGS data analysis and support the largescale identification of variants,

131    NGSpop provides a one-step user mode that can run all processes employed by NGSpop

132    with a single click. When NGSpop runs using the one-step mode, the default options will

133    be used for each step. The user can customize the default options used in the one-step

134    mode by first using the step-by-step mode.

135

## Quality control

137    To identify highly accurate genomic variation information from the population, the

138    quality of the NGS data should be carefully checked and filtered; FastQC (version 0.11.5)

139    is used for this purpose in NGSpop. Sequence reads that are below the score (Phred) [18]

140    specified by the user will be filtered out and low-quality regions at the 5′- and 3′-ends can

141    be trimmed using TrimmOmatic (version 0.36) [13].

7

142

## Read mapping and duplicate removal

144    NGSpop employs BWA (version 0.7.16a) [14] as a mapping tool for the NGS reads. To

145    convert the BWA sequence alignment map format (sam) to a binary alignment map (bam)

146    format, and then to sort and index the file, NGSpop uses SAMtools [15]. If the mate-pair

147    information is not concordant with the sample library information, it should be verified

148    and fixed. For this purpose, the Fixmate command of Picard (version 2.9.4) is used in

149    NGSpop. In addition, duplicate reads are removed using the MarkDuplicates and

150    AddOrReplaceReadGroups commands of Picard, and to calculate the statistics of the

151    sequence reads, BamTools [19] is used.

152

## SNP/INDEL identification

154    Using NGSpop, the user can select a SNP/INDEL identification algorithm from the

155    Genome Analysis Toolkit (version 3.7.0) or DeepVariant (version 0.5.1). The Genome

156    Analysis Toolkit (version 3.7.0) [9] is a standard tool for single nucleotide polymorphism

157    (SNP)/INDEL identification from NGS data. To realign the reads around the INDELs,

158    NGSpop, uses the RealignerTargetCreator and IndelRealigner commands of GATK.

159    After the realignment of the reads, UnifiedGenotyper is used as a variant caller in

160    NGSpop.

161

## DeepVariant

8

163    DeepVariant is a variant caller developed by Google Inc. The tool showed overwhelming

164    quality and imputation reference performance compared to well-established pipelines

165    such as GATK [20]. NGSpop includes DeepVariant in its pipeline, and the user can select

166    between the GATK or DeepVariant algorithms. DeepVariant is a deep learning-based

167    variant caller that uses aligned reads (in BAM or CRAM format) to produce pileup image

168    tensors, and each tensor is classified using a convolutional neural network, and finally

169    reports the results in a standard VCF or gVCF file. DeepVariant supports germline variant

170    calling in diploid organisms.

171

## 172    Variant merge

173    Vcf files that are produced using the GATK or DeepVariant algorithms in the same

174    project will be merged in NGSpop. The vcf-merge script in VCFtools [20] is employed

175    to integrate the vcf files. VCF tools are a program package of perl modules and C++

176    programs.

177

## 178    Variant annotation

179    Variations in the nucleotides can change the amino acids of the genes and thus affect the

180    organism. Therefore, the functional effects of these variants on the genes should be

181    predicted. To annotate the identified variants, NGSpop uses SnpEff (version 4.3q) [16].

182    In this study, NGSpop only included the *Arabidopsis thaliana* database (TAIR10 genome

183    [22]) for SnpEff. In other studies, the SnpEff database should be included for the

9

184  appropriate organism if it is available. If there is no available database for non-model

185  organisms, the database should be generated manually.

186

## **Variant visualization**

188  The annotated variant information can be visualized using JBrowser (version 1.12.3; Fig.

189  2) [17] in NGSpop. There are four feature tracks in the JBrowser window: reference

190  sequence, annotation information of the reference in GFF, mapped reads in bam format,

191  and annotated variants. The tracks can be shown or hidden by clicking the check box of

192  the corresponding feature tracks that the user wants to investigate. The annotated variant

193  file can be downloaded by clicking the VCF file download button on the top right of the

194  JBrowser.

195

196  **Fig 2. Visualization of variants.** Four feature tracks are listed on the left panel of the

197  JBrowser: reference sequence, annotation information of reference in GFF, mapped

198  reads, and annotated variants. Only a .vcf file can be displayed in the JBrowser when

199  multiple NGS data are selected for variant analysis after the merging of multiple .vcf

200  files.

201

202

## **Results and discussion**

204  The aim of NGSpop is to provide users with an easy-to-use environment for NGS data

205  analysis, regardless of whether the user is an expert in bioinformatics. To this end,

10

206    NGSpop provides users with two modes: a step-by-step mode for beginners and a one-

207    step mode for experts. There are currently many workflows and easy-to-use tools

208    available for NGS analysis, but as the user is required to run each step manually and wait

209    until each step ends before proceeding, they can slow the rate of analysis. Furthermore,

210    these tools were not designed for population studies, and they only provide users with a

211    step-by-step mode or a difficult hierarchical workflow design. Some tools do provide

212    user-bash script interfaces, but these can be difficult to learn. However, when using

213    NGSpop, only a single click of the run button is required and the results can be visualized

214    using JBrowser. Even though the software provides a user graphic interface, NGSpop

215    accepts multiple pairs of fastq files to support population-level studies. Thereby, users

216    can identify variants in large scale datasets from population studies using only their

217    personal computers (PCs) or workstations. Moreover, NGSpop provides a selection of

218    variant calling algorithms from GATK and DeepVariant in the variant calling steps.

219    Variant calling is an important step in NGS data analysis and genetic studies. There are

220    many tools that identify high-quality and reliable variants from NGS data, but none of

221    them can identify all variants. Therefore, researchers have used and combined multiple

222    tools to identify variants from NGS data. To assess the coverage of the variants identified,

223    we compared the variant calls from GATK and DeepVariant when using NGSpop. For

224    the benchmark test, we generated a total of five test datasets using the complete

225    *Arabidopsis thaliana* genome sequencing data from the DNA Data Bank of Japan (DDBJ)

226    FTP site under the accession number SRR519473 (paired-end run with 52,154,720 reads

227    and 10,430,944,000 bp) [24]. The sequencing data were generated by the *Arabidopsis*

228    *thaliana* 1001 genome project (http://1001genomes.org) [23] using the Illumina HiSeq

11

229    2000 platform. The detailed specifications of the benchmark system and benchmark

230    results are summarized in Table 3. NGSpop took a total of 2 h 46 m 46 s to go from the

231    raw reads to variant annotation or visualization for the five test datasets using GATK,

232    whereas it took 1 h 48 min 00 s when using DeepVariant. A total of 113,163 variants

233    were identified using GATK, and 128,530 variants were identified using DeepVariant.

234    Among the identified variants, 111,918 overlapped between GATK and DeepVariant.

235    Meanwhile, 1,245 and 16,612 were specific to GATK and DeepVariant, respectively.

236    DeepVariant with Tensorflow was faster than GATK in variant calls and identified more

237    variants than GATK. Consequently, NGSpop was found to be an easy-to-use platform for

238    variant calling using GATK and DeepVariant.

239

240    **Table 3. Comparison of the variant calling algorithms used in the NGSpop software.**

| Benchmark Machine | | | | | Variant calling algorithms | | | | |
|---|---|---|---|---|---|---|---|---|---|
| CPU | RAM (GBytes) | Storage (TBytes) | OS | | GATK | | | DeepVariant | |
| | | | | Dataset | SNP | Time | Intersection | SNP | Time |
| | | | | | Count | (hh:mm:ss) | | Counts | (hh:mm:ss) |
| Intel(R) Xeon(R) CPU E5-2609 v3 @ 1.90GHz | 32 | 1.8 | Ubuntu 18.04 | 1 | 20,896 | 0:32:15 | 20,671 | 23,761 | 0:21:19 |
| | | | | 2 | 20,404 | 0:32:24 | 20,201 | 23,207 | 0:20:04 |
| | | | | 3 | 22,952 | 0:33:02 | 22,705 | 26,087 | 0:21:38 |
| | | | | 4 | 24,616 | 0:35:11 | 24,338 | 27,920 | 0:22:35 |
| | | | | 5 | 24,295 | 0:33:54 | 24,003 | 27,555 | 0:22:24 |
| | | | | Sum | 113,163 | 2:46:46 | 111,918 | 128,530 | 1:48:00 |
| | | | | Average | 22,633 | 0:33:21 | 22,384 | 25,706 | 0:21:36 |

241    Benchmark tests were performed for NGSpop using GATK and DeepVariant as the

242    variant calling algorithms.

243

244

# Conclusions

246   Large-scale parallel sequencing has become a popular tool to identify sequence variations,

247   and many tools have now been developed to analyze NGS data. Although many tools

248   have been developed, few support population-level or cohort-level sequencing data.

249   Owing to the lack of population-level analysis tools, many researchers find it difficult to

250   analyze the massive volumes of NGS data that they produce. Researchers should ideally

251   write scripts to analyze NGS data on the Linux command line. NGSpop is a user-friendly

252   software for researchers who are not familiar with the command line interface and do not

253   want to write shell scripts. Therefore, NGSpop provides the user with an easy-to-use

254   interface and helps to automate the detection of variations from the NGS data at the

255   population level. NGSpop helps genomics researchers who want to analyze population-

256   level NGS data with an easy-to-use GUI. We developed NGSpop to support population-

257   level NGS data analysis; however, there are some limitations. First of all, NGSpop only

258   accepts FASTQ format data that has been produced using Illumina platform because there

259   are too many parameters to consider when analyzing all types of NGS platforms. Next,

260   NGSpop only supports Linux operating system because DeepVariant, one of the variant

261   calling algorithms, can only be used with Linux operating systems. In future studies, we

262   will include functionalities that support NGS platforms other than Illumina while

263   accounting for the variations in formats.

264

265   **Availability and requirements**

266   **Project name:** NGSpop

13

267 **Project home page:** https://sourceforge.net/projects/ngspop/

268 **Operating system(s):** Linux

269 **Programming language:** JavaFX

270 **Other requirements:** All Perl libraries are listed in the Supplementary information.

271 **License:** GNU General Public License

272 **Any restrictions to use by non-academics:** license needed

273

274

275

276

277 **Availability of data and materials**

278 **Test datasets:** The test datasets used as the whole-genome shotgun sequencing data are

279 available from the project home page. (https://sourceforge.net/projects/ngspop/).

280

281

282 **Acknowledgements**

283 Not applicable

284

285 # References

286 1. Ng SB, Turner EH, Robertson PD, Flygare SD, Bigham AW, Lee C, et al. Targeted
287    capture and massively parallel sequencing of 12 human exomes. Nature.
288    2009;461(7261): 272-276.
289 2. Ng SB, Buckingham KJ, Lee C, Bigham AW, Tabor HK, Dent KM, et al. Exome
290    sequencing identifies the cause of a Mendelian disorder. Nat Genet. 2010;42(1):
291    30-35.
292 3. Ng SB, Bigham AW, Buckingham KJ, Hannibal MC, McMillin MJ, Gildersleeve HI,
293    et al. Exome sequencing identifies MLL2 mutations as a cause of Kabuki
294    syndrome. Nat Genet. 2010;42(9): 790-793.

295  4. Splinter K, Adams DR, Bacino CA, Bellen HJ, Bernstein JA, Cheatle-Jarvela AM, et
296     al. Effect of genetic diagnosis on patients with previously undiagnosed disease. N
297     Engl J Med. 2018;379(22): 2131-2139.
298  5. Yu H, Xie W, Wang J, Xing Y, Xu C, Li X, et al. Gains in QTL detection using an
299     ultra-high density SNP map based on population sequencing relative to traditional
300     RFLP/SSR markers. PLOS ONE. 2011;6(3): e17595.
301  6. Schaeffer LR. Strategy for applying genome-wide selection in dairy cattle. J Anim
302     Breed Genet. 2006;123(4): 218-223.
303  7. 1000 Genomes Project Consortium, Abecasis GR, Altshuler D, Auton A, Brooks LD,
304     Durbin RM, et al. A map of human genome variation from population-scale
305     sequencing. Nature. 2010;467(7319): 1061-1073
306  8. Blanca JM, Pascual L, Ziarsolo P, Nuez F, Cañizares J. ngs_backbone: a pipeline for
307     read cleaning, mapping and SNP calling using next generation sequence. BMC
308     Genomics. 2011;12: 285.
309  9. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
310     genome analysis toolkit: a MapReduce framework for analyzing next-generation
311     DNA sequencing data. Genome Res. 2010;20(9): 1297-1303.
312  10. Goecks J, Nekrutenko A, Taylor J. Galaxy Team. Galaxy: a comprehensive approach
313     for supporting accessible, reproducible, and transparent computational research in
314     the life sciences. Genome Biol. 2010;11(8): R86.
315  11. Liu CH, Di YP. Analysis of RNA sequencing data using CLC genomics workbench.
316     Methods Mol Biol. 2020;2102: 61-113.
317  12. Ryan Poplin, Pi-Chuan Chang, David Alexander, Scott Schwartz, Thomas Colthurst,
318     Alexander Ku, et al. A universal SNP and small-indel variant caller using deep
319     neural networks. Nat Biotechnol. 2018;36(10): 983-987.
320  13. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina
321     sequence data. Bioinformatics. 2014;30(15): 2114-2120.
322  14. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
323     transform. Bioinformatics. 2009;25(14): 1754-1760.
324  15. Li H. A statistical framework for SNP calling, mutation discovery, association
325     mapping and population genetical parameter estimation from sequencing data.
326     Bioinformatics. 2011;27(21): 2987-2993.
327  16. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for
328     annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:
329     SNPs in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3. Fly.
330     2012;6(2): 80-92.
331  17. Skinner ME, Uzilov AV, Stein LD, Mungall CJ, Holmes IH. JBrowse: a next-
332     generation genome browser. Genome Res. 2009;19(9): 1630-1638.
333  18. Nowrousian M. Next-generation sequencing techniques for eukaryotic
334     microorganisms: sequencing-based solutions to biological problems. Eukaryot
335     Cell. 2010;9(9): 1300-1310.
336  19. Ewing B, Green P. Base-calling of automated sequencer traces using Phred. II. Error
337     probabilities. Genome Res. 1998;8(3): 186-194.
338  20. Yun T, Li H, Chang PC, Lin MF, Carroll A, McLean CY. Accurate, scalable cohort
339     variant calls using DeepVariant and GLnexus. Bioinformatics 2021;36(24): 5582-
340     5589.

15

341  21. Barnett DW, Garrison EK, Quinlan AR, Strömberg MP, Marth GT. BamTools: a C++
342      API and toolkit for analyzing and managing BAM files. Bioinformatics.
343      2011;27(12): 1691-1692.
344  22. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The
345      variant call format and VCFtools. Bioinformatics 2011;27(15): 2156-2158.
346  23. Swarbreck D, Wilks C, Lamesch P, Berardini TZ, Garcia-Hernandez M, Foerster H,
347      et al. The Arabidopsis Information Resource (TAIR): gene structure and function
348      annotation. Nucleic Acids Res. 2008;36: D1009-D1014.
349  24. Long Q, Rabanal FA, Meng D, Huber CD, Farlow A, Platzer A, et al. Massive
350      genomic variation and strong selection in Arabidopsis thaliana lines from Sweden.
351      Nat Genet. 2013;45(8): 884-890.

352

353

354

355

356

357

358

# Supporting information

360  **Additional file 1. Supplementary.docx**
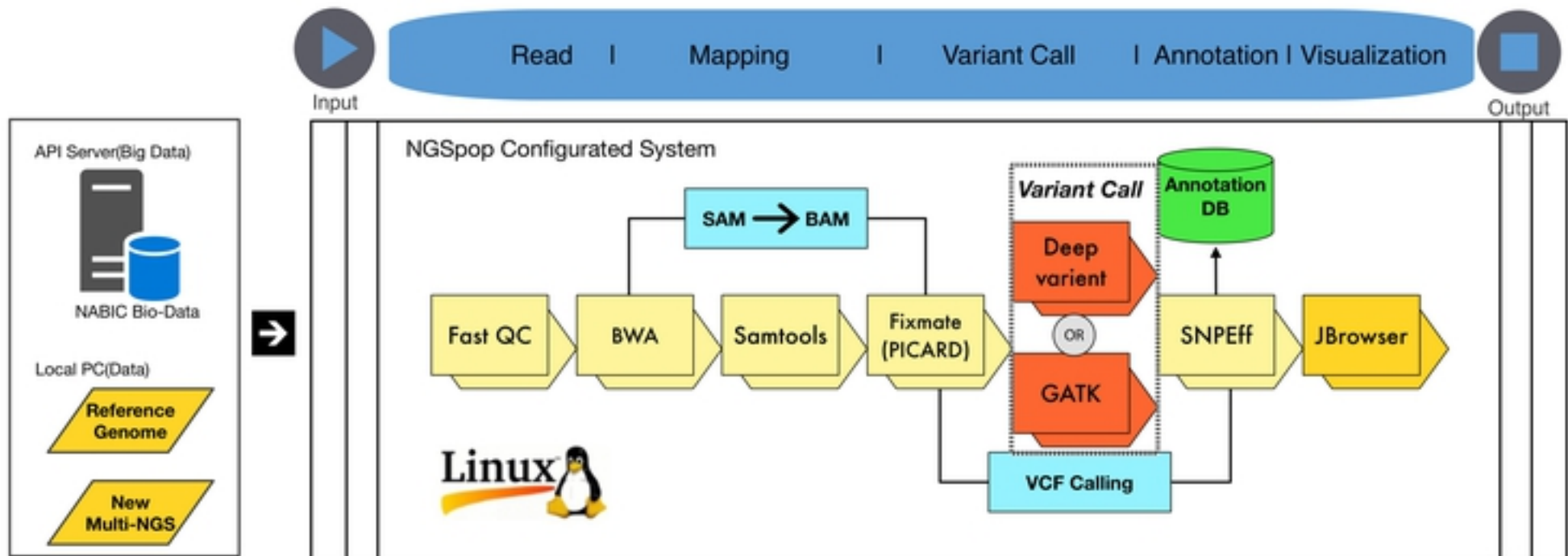
361

Figure1

Figure2